

1
2
3
4

From Univariate to Multiple Linear Regression

Univariate Linear Regression Recap:

- Model:
 $f_w, b(x) = wx + b$
 - One feature (e.g., house size) used to predict the target (e.g., price).
-

Introduction to Multiple Features

New Scenario:

We now consider **multiple features** to predict housing price:

- x_1 : Size (sq. ft)
- x_2 : Number of bedrooms
- x_3 : Number of floors
- x_4 : Age of house (in years)

This gives us **more information** for making better predictions.

A
B
C
D

New Notations

Term	Meaning
x_j	The j-th feature (e.g., size, bedrooms, etc.)
n	Total number of features
$x^{(i)}$	Feature vector of the i-th training example (a row vector)
$x_j^{(i)}$	Value of the j-th feature for the i-th example
$\vec{x}^{(i)}$	Vector notation to emphasize it's a list, not a scalar

E
F
G
H

Example:

- $x^{(2)} = [1416, 3, 2, 40]$
- $x_3^{(2)} = 2$ (3rd feature = 2 floors)



Multiple Linear Regression Model

Expanded Form:

$$f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

Example Model (housing price in \$1000s):

$$f_{w,b}(x) = 0.1x_1 + 4x_2 + 10x_3 - 2x_4 + 80$$

Interpretation of Coefficients:

- **b** = **80**: base price is \$80,000
- **0.1**: +\$100 per extra square foot
- **4**: +\$4000 per bedroom
- **10**: +\$10,000 per floor
- **-2**: -\$2000 per year older



General Form with **n** Features

$$f_{w,b}(x) = \sum_{j=1}^n w_j x_j + b$$



Vector Notation

To make the math cleaner, use vectors:

- Let:

$$\vec{w} = [w_1, w_2, \dots, w_n] \quad (\text{parameter vector})$$

$$\vec{x} = [x_1, x_2, \dots, x_n] \quad (\text{feature vector})$$

- b : remains a scalar

🔗 Dot Product (Inner Product)

🔗 Dot Product (Inner Product)

$$\vec{w} \cdot \vec{x} = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

Final Model (Compact Form):

$$f_{w,b}(x) = \vec{w} \cdot \vec{x} + b$$

⚠️ Dot product only works when both w and x are of the **same length** (same number of features).

📌 Terminology

Term	Meaning
Univariate Linear Regression	Regression with a single feature
Multiple Linear Regression	Regression with multiple features
Multivariate Regression	A different concept (not covered here); often refers to predicting multiple outputs , not using multiple inputs

Question

In the training set below, what is $x_1^{(4)}$? Please type in the number below (this is an integer such as 123, no decimal points).

Size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price (\$) in \$1000's
x_1	x_2	x_3	x_4	$j=1 \dots 4$
2104	5	1	45	$n=4$
1416	3	(2)	40	
1534	3	2	30	
852	2	1	36	
...

$x_j = j^{\text{th}}$ feature
 $n = \text{number of features}$
 $\vec{x}^{(i)} = \text{features of } i^{\text{th}} \text{ training example}$
 $x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example}$

$\vec{x}^{(2)} = [1416 \ 3 \ (2) \ 40]$
 $x_3^{(2)} = 2$

Skip Submit

9:27 / 9:51 1.25x

852 ans

🧠 What is Vectorization?

Vectorization is a method of rewriting code (especially loops and repetitive operations) to leverage **matrix and vector operations**.

📈 Why Vectorization Matters in Machine Learning?

- **Shorter Code** → Easier to read, write, and maintain.
- **Faster Execution** → Leverages modern linear algebra libraries and hardware acceleration (CPU/GPU).
- Used extensively in machine learning for **efficient computation**.

🧪 Example: Without and With Vectorization

Setup:

Let's say:

- $\vec{w} = [w_1, w_2, w_3]$
- $\vec{x} = [x_1, x_2, x_3]$
- b is a scalar
- $n = 3$

♦ **Implementation without vectorization (manual multiplication):**

```
python
CopyEdit
f = w[0] * x[0] + w[1] * x[1] + w[2] * x[2] + b
```

👎 Not scalable: what if $n = 100,000$?

♦ **Loop-based (still not vectorized):**

```
python
CopyEdit
f = 0
for j in range(n):    # j = 0 to n-1
    f += w[j] * x[j]
f += b
```

⚠ Still sequential and not fast when n is very large.

✓ **Vectorized Implementation (Best Practice):**

✨ **Math:**

$$f_{\vec{w}, b}(x) = \vec{w} \cdot \vec{x} + b$$

✨ **Code (NumPy):**

```
import numpy as np
```

```
f = np.dot(w, x) + b
```

- `np.dot(w, x)` performs the **dot product**.
 - Only **one line** does the full computation.
-

⚡ Why is Vectorization So Fast?

1. Parallel Execution

- NumPy uses optimized **C/C++ backends**.
- Operations are **executed in parallel** across CPU or **GPU** cores.

2. No Python Overhead

- For loops in Python are slow due to interpreter overhead.
 - Vectorized code executes in **compiled code** under the hood.
-



Recap: Benefits of Vectorization

Benefit	Description
Shorter Code	One-liner instead of multiple lines
Faster Code	Uses efficient C/Fortran implementations and parallelism
Scalable	Handles large <code>n</code> effortlessly
Clearer	Easier to read and debug

Parameters and features

 $\vec{w} = [w_1 \ w_2 \ w_3] \quad n=3$

b is a number

 $\vec{x} = [x_1 \ x_2 \ x_3]$

linear algebra: count from 1

```
w[0] w[1] w[2]
```

```
w = np.array([1.0, 2.5, -3.3])
b = 4
x = np.array([10, 20, 30])
```

code: count from 0

Without vectorization $n=100,000$

 $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

```
f = w[0] * x[0] +
    w[1] * x[1] +
    w[2] * x[2] + b
```



Without vectorization

$$f_{\vec{w}, b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \rightarrow j=1 \dots n$$

$\text{range}(0, n) \rightarrow j=0 \dots n-1$

```
f = 0
for j in range(0, n):
    f = f + w[j] * x[j]
f = f + b
```



Vectorization

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

```
f = np.dot(w, x) + b
```

⚙️ How Vectorization Works (Behind the Scenes)

🧠 Sequential vs Parallel Computation

🔁 Without Vectorization (For Loop)

- Executes operations **step-by-step, one after another.**

```
for j in range(16):
    result[j] = w[j] * x[j]
```

- Time steps:

- t0t_0t0: compute $w_0 \times x_0 w_0 \times x_0$
- t1t_1t1: compute $w_1 \times x_1 w_1 \times x_1$
- ...
- t15t_{15}t15: compute $w_{15} \times x_{15} w_{15} \times x_{15}$

▼ **Slow for large feature sets** (like $n = 10,000$).

⚡ With Vectorization (Parallel Processing)

```
\  
result = np.dot(w, x)
```

- Uses **specialized hardware** (CPU/GPU)
- Computes all 16 multiplications **simultaneously**
- Then adds the 16 results efficiently with **parallel summation**

▲ **Faster**, especially for large **n**.



Use Case: Parameter Updates in Gradient Descent



Goal:

Update weights w_j using gradient descent:

$$w_j := w_j - \alpha \cdot d_j$$

Where:

- α : learning rate (e.g., 0.1)
 - d_j : derivative of cost w.r.t w_j
-

✗ Without Vectorization

```
for j in range(16):  
    w[j] = w[j] - 0.1 * d[j]
```

- Updates each weight **individually** (slower)
-

✓ With Vectorization

```
w = w - 0.1 * d
```

- All 16 updates done in one line
 - Leverages NumPy array broadcasting
 - Runs faster using parallel computation
-

Practical Example in NumPy

```
import numpy as np

# Parameters
w = np.array([w1, w2, ..., w16])
d = np.array([d1, d2, ..., d16])
alpha = 0.1

# Vectorized update
w = w - alpha * d
```

Real Impact

Situation	Speed Gain
Small n (e.g., 16 features)	Slight
Large n (e.g., 10,000+ features)	Huge
Large training dataset (m large)	Critical

 Vectorization can make the **difference between seconds and hours** in ML training!

Without vectorization

```
for j in range(0,16):
    f = f + w[j] * x[j]
```

t_0

t_1

\dots

t_{15}

$f + w[0] * x[0]$

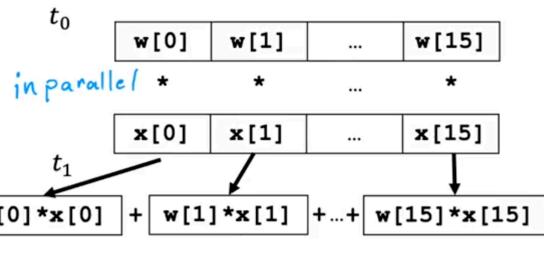
$f + w[1] * x[1]$

\dots

$f + w[15] * x[15]$

Vectorization

```
np.dot(w, x)
```



Gradient descent $\vec{w} = (w_1 \quad w_2 \quad \dots \quad w_{16})$ ~~b~~ parameters

derivatives $\vec{d} = (d_1 \quad d_2 \quad \dots \quad d_{16})$

```
w = np.array([0.5, 1.3, ... 3.4])
d = np.array([0.3, 0.2, ... 0.4])
```

compute $w_j = w_j - 0.1d_j$ for $j = 1 \dots 16$

Without vectorization

```
w1 = w1 - 0.1d1
w2 = w2 - 0.1d2
:
w16 = w16 - 0.1d16
```

```
for j in range(0,16):
    w[j] = w[j] - 0.1 * d[j]
```

With vectorization

$$\vec{w} = \vec{w} - 0.1\vec{d}$$

```
w = w - 0.1 * d
```



Gradient Descent with Multiple Features (Vectorized Form)



Recap: Multiple Linear Regression in Vector Form

- Model:

$$f_{\vec{w}, b}(x) = \vec{w} \cdot \vec{x} + b$$

- $\vec{w} = [w_1, w_2, \dots, w_n]$: parameter vector (length = number of features)
- $\vec{x} = [x_1, x_2, \dots, x_n]$: feature vector
- b : scalar bias term

⌚ Cost Function (Vectorized Notation)

- Let $J(\vec{w}, b)$ be the cost function (Mean Squared Error):

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2$$

- m : number of training examples
- Each $x^{(i)}$ is a **vector** of features



Gradient Descent: Updating Parameters

📈 For One Feature:

- Update rules:

$$\begin{aligned} w &:= w - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \\ b &:= b - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)}) \end{aligned}$$

📊 For Multiple Features:

- For each w_j (from 1 to n):

$$w_j := w_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

- Update b as before:

$$b := b - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})$$



⚡ Vectorized Gradient Descent in Code (using NumPy)

```
# X: shape (m, n) -> each row is a training example
# y: shape (m, ) -> actual outputs
```

```

# w: shape (n, ) -> parameters
# b: scalar
# alpha: learning rate

# Prediction
f = np.dot(X, w) + b                      # shape (m, )

# Error term
error = f - y                             # shape (m, )

# Gradients
dj_dw = (1/m) * np.dot(X.T, error)    # shape (n, )
dj_db = (1/m) * np.sum(error)          # scalar

# Update
w = w - alpha * dj_dw
b = b - alpha * dj_db

```

 Everything is done **efficiently in parallel**, thanks to vectorization.

Summary of Vectorized Gradient Descent

Component	Description
<code>np.dot(X, w)</code>	Matrix-vector multiplication: all predictions at once
<code>np.dot(X.T, error)</code>	Efficient computation of all $w_j w_j$ gradients
Vectorization	Speeds up large-scale training and simplifies implementation

Aside: The Normal Equation

What is it?

- An **analytical** solution to linear regression:

$$\vec{w} = (X^T X)^{-1} X^T y$$

Characteristics:

Feature	Description
✓ No iterations	Directly solves for $\vec{w} \backslash \text{vec}\{w\} w$ and $b b b$
✗ Slow for large n	Matrix inversion is computationally expensive
✗ Not generalizable	Only works for linear regression, not for logistic regression, neural networks, etc.

Most practitioners don't use it directly, but some ML libraries may use it under the hood.

	Previous notation	Vector notation
Parameters	w_1, \dots, w_n b	\vec{w} <i>vector of length n</i> b <i>still a number</i>
Model	$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(\underbrace{w_1, \dots, w_n, b})$	$J(\vec{w}, b)$ <i>dot product</i>

Gradient descent

<pre>repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{w_1, \dots, w_n, b})$ $b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n, b})$ }</pre>	<pre>repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$ $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$ }</pre>
---	---

Gradient descent

<p>One feature</p> <pre> repeat { $w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$ ↳ $\frac{\partial}{\partial w} J(w, b)$ $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$ simultaneously update w, b } </pre>	<p>n features ($n \geq 2$)</p> <pre> repeat { $j=1$ $w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)}$: $j=n$ $w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$ $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})$ simultaneously update w_j (for $j = 1, \dots, n$) and b } </pre>
--	--

1 / 1 point

2. Which of the following are potential benefits of vectorization? Please choose the best option.

- It makes your code run faster
- It can make your code shorter
- It allows your code to run more easily on parallel compute hardware
- All of the above

Correct

Correct! All of these are benefits of vectorization!

1 / 1 point

3. True/False? To make gradient descent converge about twice as fast, a technique that almost always works is to double the learning rate α .

- False
- True

Correct

Doubling the learning rate may result in a learning rate that is too large, and cause gradient descent to fail to find the optimal values for the parameters w and b .

🔍 Problem Setup

You're trying to **predict house prices** using two features:

- **$x_1 = \text{size of the house (sq. feet)}$** → ranges from **300 to 2000**
- **$x_2 = \text{number of bedrooms}$** → ranges from **0 to 5**

That's a **big range difference** between x_1 and x_2 .

Why Is This a Problem?

If features have very **different scales**, then:

- **Large feature values** (like size in sq. feet) will **dominate the cost function's sensitivity**.
- This leads to **elongated contours** (think of stretched ovals) in the cost function's landscape.

As a result:

- Gradient descent will take **zig-zag paths**, bouncing slowly down the valley.
 - It takes **much longer to converge** to the optimal values of weights (w_1, w_2).
-

Example to Understand

Let's say you have a house:

- Size: **2000 sq. feet**
- Bedrooms: **5**
- Price: **\$500,000 (500k)**

Case 1: Parameters → $w_1 = 50, w_2 = 0.1, b = 50$

Predicted price:

```
yaml
CopyEdit
= 50 × 2000 + 0.1 × 5 + 50
= 100,000 + 0.5 + 50 ≈ 100,100.5k → Way too high!
```

→ **Bad prediction.**

Case 2: Parameters → $w_1 = 0.1, w_2 = 50, b = 50$

Predicted price:

$$\begin{aligned} &= 0.1 \times 2000 + 50 \times 5 + 50 \\ &= 200 + 250 + 50 = 500k \rightarrow \text{Perfect!} \end{aligned}$$

→ Good prediction.

Observation:

- Large-valued features (like x_1) → **Small weights**
 - Small-valued features (like x_2) → **Large weights**
-



Cost Function Visualized

- x-axis = w_1 (weight for size), y-axis = w_2 (weight for bedrooms)
- Because of the unbalanced feature ranges:
 - The **cost function contour is stretched** (elliptical).
 - Small steps in w_1 make **big prediction changes**.
 - Steps in w_2 hardly change predictions at all.

So gradient descent **struggles to move efficiently**.



Solution: Feature Scaling

✓ What is Feature Scaling?

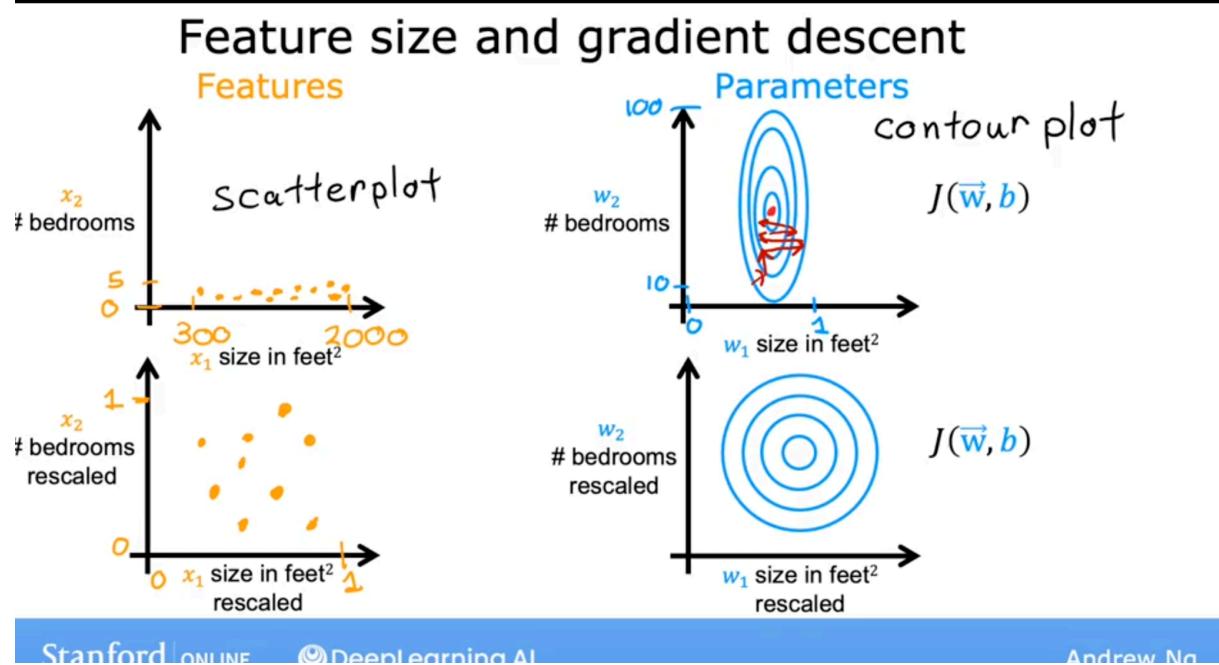
Rescale features so they all take **similar ranges** (e.g., from 0 to 1).

Example:

Feature	Original Range	After Scaling
x_1 (size)	300 to 2000	0 to 1
x_2 (bedrooms)	0 to 5	0 to 1

Now:

- Cost function contours are **circular**.
- Gradient descent moves **directly toward the minimum**.
- **Convergence is faster and more stable.**



Feature Scaling: How to Implement It

When your features (input variables) have very different **ranges**, gradient descent can become inefficient. Feature scaling solves this.



Method 1: Min-Max Scaling (Rescaling to [0,1])

 **Formula:**

$$x_{\text{scaled}} = \frac{x}{x_{\text{max}}}$$

 **Example:**

- Feature x_1 = house size $\in [300, 2000]$:

$$x_1^{\text{scaled}} = \frac{x_1}{2000} \Rightarrow x_1^{\text{scaled}} \in [0.15, 1]$$

- Feature x_2 = number of bedrooms $\in [0, 5]$:

$$x_2^{\text{scaled}} = \frac{x_2}{5} \Rightarrow x_2^{\text{scaled}} \in [0, 1]$$

Method 2: Mean Normalization (Center around 0)

 **Formula:**

$$x_{\text{normalized}} = \frac{x - \mu}{x_{\text{max}} - x_{\text{min}}}$$

- μ : mean of the feature
- $x_{\text{max}}, x_{\text{min}}$: max and min values in dataset

 **Example:**

- Suppose:

- $\mu_1 = 600, x_1 \in [300, 2000]$

$$x_1^{\text{normalized}} = \frac{x_1 - 600}{2000 - 300} \Rightarrow x_1^{\text{normalized}} \in [-0.18, 0.82]$$

- For x_2 , if $\mu_2 = 2.3, x_2 \in [0, 5]$:

$$x_2^{\text{normalized}} = \frac{x_2 - 2.3}{5} \Rightarrow x_2^{\text{normalized}} \in [-0.46, 0.54]$$



Method 3: Z-Score Normalization (Standardization)

12
34 **Formula:**

$$x_{\text{z-score}} = \frac{x - \mu}{\sigma}$$

- μ : mean
- σ : standard deviation

✓ **Example:**

- Feature 1: $\mu = 600, \sigma = 450$

$$x_1^z = \frac{x_1 - 600}{450} \Rightarrow x_1^z \in [-0.67, 3.1]$$

- Feature 2: $\mu = 2.3, \sigma = 1.4$

$$x_2^z = \frac{x_2 - 2.3}{1.4} \Rightarrow x_2^z \in [-1.6, 1.9]$$

✓ Rule of Thumb

You ideally want **all features to lie between**:

- **-1 to +1**
- Or loosely within **[-3, +3]** is acceptable
- Avoid ranges like **[0, 100]** or **[-100, 100]**



When to Scale?

Feature Range	Scale ?	Why?
[0.001, 0.01]	✓	Too small → weak gradient signal
[-100, 100]	✓	Too large → zig-zag descent
[-1, 1]	✗	Already in ideal range
[0, 3]	✗	Acceptable
[98.6, 105] (body temp in °F)	✓	Around 100 → can slow down GD

Want Code for It?

Here's a quick example in Python using NumPy:

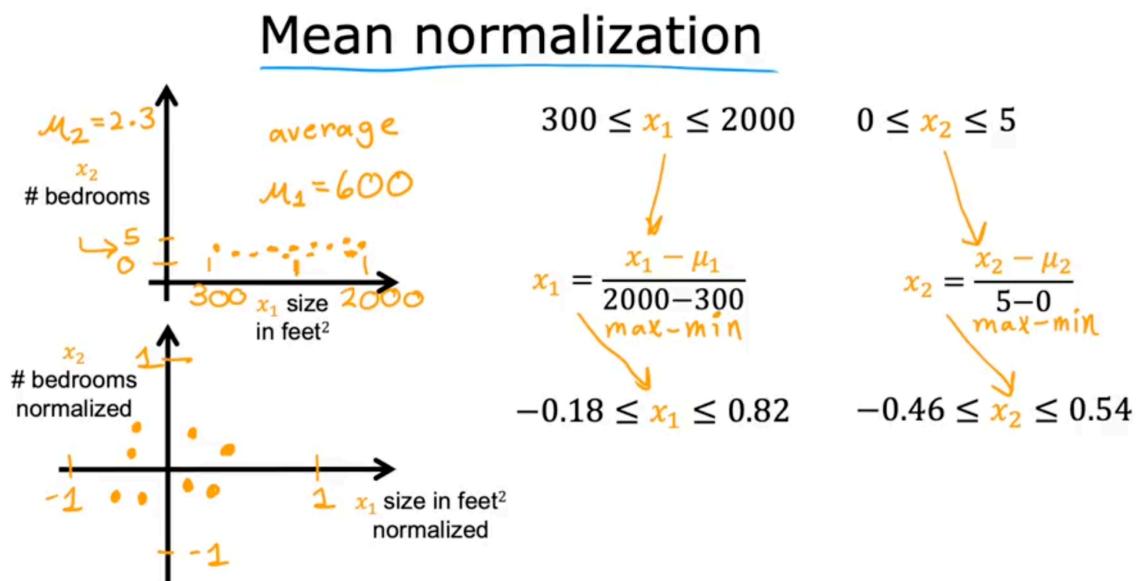
```
import numpy as np

# Example feature
x1 = np.array([300, 600, 1000, 1500, 2000])

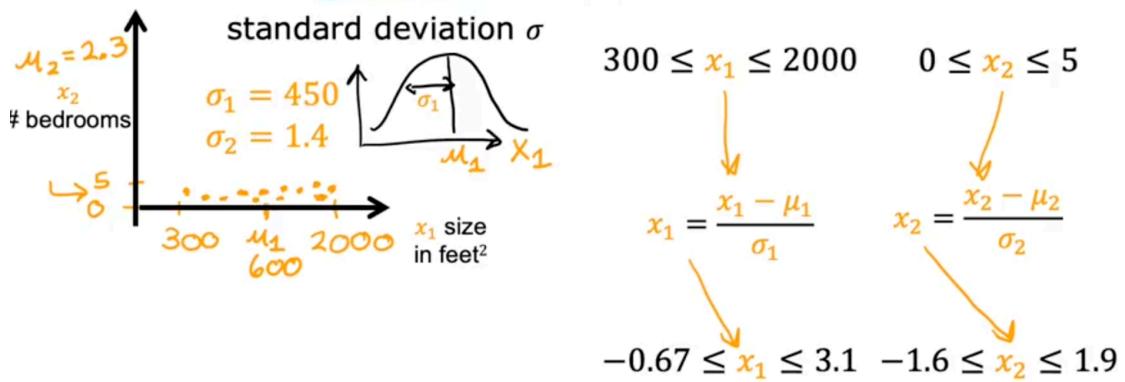
# 1. Min-Max Scaling
x1_minmax = x1 / np.max(x1)

# 2. Mean Normalization
mu1 = np.mean(x1)
x1_meannorm = (x1 - mu1) / (np.max(x1) - np.min(x1))

# 3. Z-score Normalization
sigma1 = np.std(x1)
x1_zscore = (x1 - mu1) / sigma1
```



Z-score normalization



Feature scaling

aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$-3 \leq x_j \leq 3$	}	acceptable ranges
$-0.3 \leq x_j \leq 0.3$		

$0 \leq x_1 \leq 3$	Okay, no rescaling
$-2 \leq x_2 \leq 0.5$	Okay, no rescaling
$-100 \leq x_3 \leq 100$	too large → rescale
$-0.001 \leq x_4 \leq 0.001$	too small → rescale
$98.6 \leq x_5 \leq 105$	too large → rescale

How to Know If Gradient Descent Is Converging?

The **goal of gradient descent** is to find parameters w and b that **minimize the cost function** $J(w, b)$. So we want to monitor if J is **decreasing steadily** over time.

Step 1: Plot the Learning Curve

What is a learning curve?

- **X-axis:** Number of gradient descent iterations
- **Y-axis:** Value of cost function $J(w, b)$

What should it look like?

- The curve should **decrease after every iteration**
- It should **flatten out** over time (that's convergence)

Interpretation:

Observation	Meaning
J decreases steadily	Gradient descent is working well 
J increases at some iteration	Learning rate might be too large  or code has a bug 
J flattens out	Gradient descent has converged 
J decreases too slowly	Learning rate might be too small 

Step 2: Use Epsilon (ε) Convergence Test

Let's define:

$$\varepsilon = \text{a small threshold, e.g., } 0.001 \text{ or } 10^{-3}$$

If the **change in cost function** between two iterations is smaller than ε :

$$|J^{(i)} - J^{(i+1)}| < \varepsilon$$

You can consider that **convergence has happened**.

Caution:

Choosing the right value for epsilon can be **tricky**. If it's too small, convergence may never be detected. That's why:

 Visualizing the learning curve is often more reliable than automatic convergence checks.

How Many Iterations Are Enough?

There's **no fixed number** — it depends on:

- The learning rate α

- The complexity of data
- The cost function shape

 Sometimes 30 iterations are enough. Other times, you may need **100,000+**.

So always **plot the learning curve** — that's your best guide.

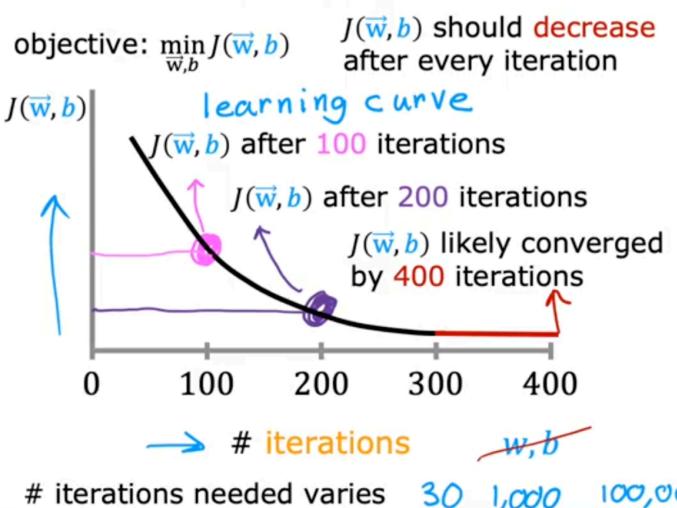


Example of a Good Learning Curve

Iteration	Cost $J(w, b)$
0	120
10	75
50	31
100	17
150	9
200	4.2
250	3.9
300	3.8 (almost flat now)

This shows **fast convergence** and then **flattening out** — ideal!

Make sure gradient descent is working correctly



Automatic convergence test
Let ϵ "epsilon" be 10^{-3} .
 0.001
If $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration,
declare convergence.
(found parameters \vec{w}, b
to get close to
global minimum)

Choosing a Good Learning Rate (α)

Goal of Learning Rate

The learning rate α controls how big each step is when updating your model parameters w and b during gradient descent:

$$w := w - \alpha \cdot \partial J / \partial w = w - \alpha \cdot \nabla_w J \\ b := b - \alpha \cdot \partial J / \partial b = b - \alpha \cdot \nabla_b J$$

What Happens If α is Too Large?

- Updates **overshoot the minimum**
- Cost function J may **bounce up and down**
- Gradient descent may **diverge**, i.e., J increases after each step
- May look like this on the learning curve:

makefile

CopyEdit

$J: 50 \rightarrow 70 \rightarrow 40 \rightarrow 80 \rightarrow 30 \rightarrow \star$ (oscillates wildly or explodes)

Common Mistake in Code:

Wrong update sign:

```
# Wrong ✗  
w = w + a * gradient
```

Should be:

```
# Correct ✓  
w = w - a * gradient
```

The minus sign is **critical** because it moves you **toward the minimum**.

What Happens If α is Too Small?

- Gradient descent works... but **very slowly**
- Takes **forever to converge**
- Example learning curve:

$J: 50 \rightarrow 49.9 \rightarrow 49.8 \rightarrow 49.7 \rightarrow \dots$ (tiny progress)

What Should a Good Learning Curve Look Like?

- Cost J decreases rapidly at first, then flattens
- No upward spikes
- Example:

$J: 50 \rightarrow 32 \rightarrow 21 \rightarrow 15 \rightarrow 8 \rightarrow 4 \rightarrow 3.8 \rightarrow 3.7 \dots$

Debug Tip (Powerful!)

Try a very small α (like 0.0001) and see if J decreases every iteration.

If it still increases, you likely have a bug in your gradient descent code.

Strategy for Choosing α

1. Try a range of α values:

Start with small values and **increase by a factor of ~3**:

0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1.00.001, \ 0.003, \ 0.01, \ 0.03, \ 0.1, \ 0.3, \ 1.00.001,
0.003, 0.01, 0.03, 0.1, 0.3, 1.0

2. For each α :

- Run **just a few iterations** (say 100)
- Plot cost J vs iterations
- Pick the **largest α** that **still causes smooth decrease in J**

You want to pick **the fastest learning rate that still behaves well**.

Recap: How to Choose α

Step	Action
p	
✓ 1	Try values like 0.001, 0.003, 0.01, 0.03, 0.1
✓ 2	Plot cost J for each α
✓ 3	If J goes up or oscillates, α is too large
✓ 4	If J decreases very slowly, α is too small
✓ 5	Choose the largest α that still gives a smooth descent



Real-World Practice

In code, your process might look like:

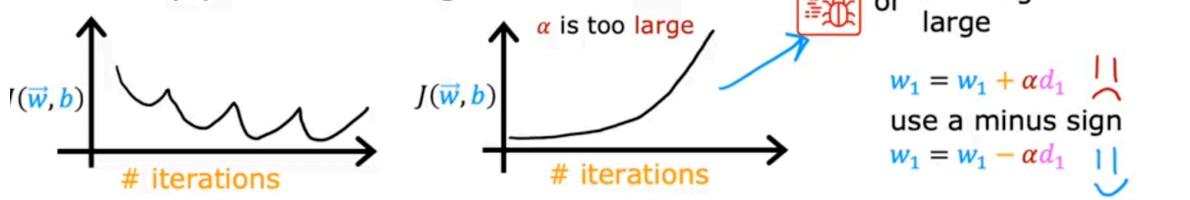
```
alphas = [0.001, 0.003, 0.01, 0.03, 0.1]
for alpha in alphas:
    print(f"Testing alpha = {alpha}")
    w, b = initialize_parameters()
    J_history = gradient_descent(X, y, w, b, alpha, iterations=100)
    plot(J_history)
```



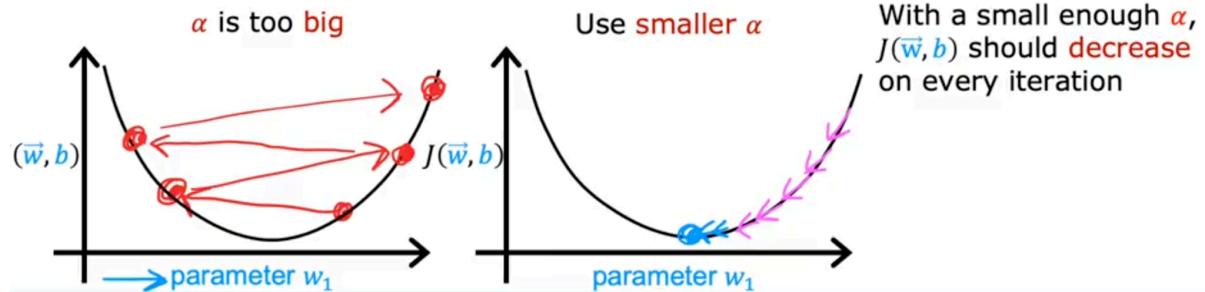
Final Insight

Tuning learning rate is a mix of science and intuition. Don't expect to find the perfect α instantly. Use plots, play around, and trust your eyes.

Identify problem with gradient descent

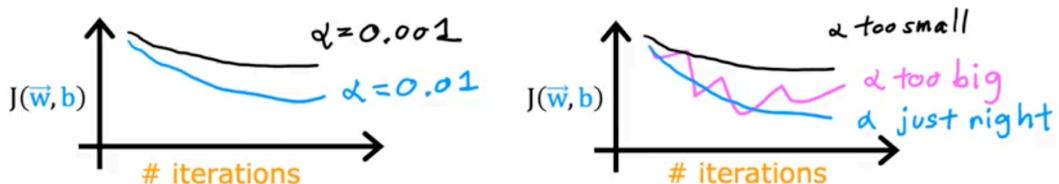


Adjust learning rate



Values of α to try:

$$\dots 0.001 \xrightarrow{3X} 0.003 \xrightarrow{\approx 3X} 0.01 \xrightarrow{3X} 0.03 \xrightarrow{\approx 3X} 0.1 \xrightarrow{3X} 0.3 \xrightarrow{\approx 3X} 1 \dots$$



You run gradient descent for 15 iterations with $\alpha = 0.3$ and compute $J(w)$ after each iteration. You find that the value of $J(w)$ increases over time. How do you think you should adjust the learning rate α ?

- Try running it for only 10 iterations so $J(w)$ doesn't increase as much.
- Keep running it for additional iterations
- Try a smaller value of α (say $\alpha = 0.1$).
- Try a larger value of α (say $\alpha = 1.0$).

Correct

Since the cost function is increasing, we know that gradient descent is diverging, so we need a lower learning rate.

What is Feature Engineering?

Feature Engineering is the process of:

- **Creating new features** from the raw data you already have.
 - **Transforming or combining features** in a way that makes your model perform better.
 - Using **intuition or domain knowledge** to guide how you represent input data to the algorithm.
-

Example: Predicting House Prices

You have 2 original features:

- x_1 : **Frontage** – width of the lot (in feet or meters)
- x_2 : **Depth** – length of the lot

A basic linear model would be:

$$f(x) = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

This uses **frontage and depth as independent predictors**.

But here's the smart observation:

 **Area = Frontage × Depth = $x_1 \times x_2$**

So we engineer a **new feature**:

- $x_3 = x_1 \cdot x_2 \rightarrow \text{Area of the land}$

Now, you build a better model:

$$f(x) = w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + b$$

The model will **learn** whether width, depth, or area (or some combination) is most predictive of house price.

Why Feature Engineering Matters

- It gives the algorithm better signals.
 - It helps capture **non-linear patterns** without needing fancy models.
 - It can **improve accuracy significantly**, especially when raw features are not sufficient.
-

Key Insight

Rather than blindly using the features given to you, think about the **real-world relationships** between them. You might be able to **combine or transform** them into more powerful predictors.

Feature engineering

$$f_{\vec{w}, b}(\vec{x}) = w_1 \underline{x_1} + w_2 \underline{x_2} + b$$

frontage depth

area = frontage × depth

*$x_3 = x_1 x_2$
new feature*

$$f_{\vec{w}, b}(\vec{x}) = \underline{w_1} x_1 + \underline{w_2} x_2 + \underline{w_3} x_3 + b$$



Feature engineering:
Using **intuition** to design
new features, by
transforming or **combining**
original features.

If you have measurements for the dimensions of a swimming pool (length, width, height), which of the following two would be a more useful engineered feature?

- length × width × height*
- length + width + height*

 **Correct**

The volume of the swimming pool could be a useful feature to use. This is the more useful engineered feature of the two.

Why Polynomial Regression?

Sometimes, your data can't be captured well with a **straight line** (like in linear regression). For example:

You have data where x is house size in square feet, and y is price.

Plotting this may show a **curve**, not a straight line. So, what do we do?

⟳ From Linear to Polynomial Regression

We **expand the feature set** by including powers of x , like:

- xxx
- $x^2x^2x^2$
- $x^3x^3x^3$
- ... up to $x^d x^d \dots x^d$ (for degree- d polynomial)

Example:

Example:

Quadratic Regression:

$$f(x) = w_1 \cdot x + w_2 \cdot x^2 + b$$

This fits a **U-shaped curve**. But in housing prices, this might not make sense, because it eventually **comes back down**, which implies that very large houses become **cheaper**, which is unrealistic.

▲ Cubic Regression:

$$f(x) = w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3 + b$$

Now you get a **flexible curve** that can better match increasing prices as house size increases.

⚠ The Importance of Feature Scaling

As you raise x to higher powers:

- If x ranges from 1 to 1,000

- Then x^2 ranges from 1 to **1,000,000**
- And x^3 ranges from 1 to **1,000,000,000**

These huge differences in scale can **slow down or break gradient descent**. That's why **feature scaling (e.g., normalization)** is essential.

Other Feature Choices

Polynomial regression doesn't just mean powers like x^2 or x^3 . You can get creative:

- Use \sqrt{x} (square root of x)
- Use $\log(x)$
- Combine features in new ways

Example:

$$f(x) = w_1 \cdot x + w_2 \cdot \sqrt{x} + b$$

The square root curve **increases more slowly** as x grows and **never comes down**, which might be a better fit for real-world housing price trends.



Feature Engineering + Polynomial Regression = Power

You're **engineering features** (transforming x into x^2, x^3, \sqrt{x} , etc.) to help the model **understand the pattern better**.

Instead of just:

$$f(x) = w_1 \cdot x + b$$

You now have:

$$f(x) = w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3 + b$$

Or:

$$f(x) = w_1 \cdot x + w_2 \cdot \sqrt{x} + b$$

The model can now fit **non-linear relationships**.

Polynomial regression

